Fast Inverse Square Roots

(0x5f3759df)

floating point $f_X =$ 

| S | $E_X$ | $M_X$ |
|---|-------|-------|

$\underset{1\,bit}{} \quad \underset{b\;bits}{} \qquad \underset{N-b-1\;bits}{}$

$\underset{N\text{ bits}}{}$

Sqrt only defined for positive values so $S = 0$

$\therefore f_X = \left(1 + \dfrac{M_X}{2^{N-b-1}}\right) 2^{E_X - (2^{b-1}-1)}$

For simplicity, let $L = 2^{N-b-1}$, which normalizes the mantissa $M$ to $0 \le M < 1$ + let $B = 2^{b-1}-1$, which is the exponent bias.

$$f_X = \left(1 + \dfrac{M_X}{L}\right) 2^{E_X - B} \qquad (1)$$

Looking for $y = \dfrac{1}{\sqrt{x}} = x^{-\frac{1}{2}}$

Let $f_X$ & $f_y$ be floating point representations of $x$ + $y$ respectively.

$\therefore f_y = f_X^{-\frac{1}{2}}$ ignoring errors introduced by floating pt. approx.

$\log_2 f_y = \log_2 \left(f_X^{-\frac{1}{2}}\right)$

$\log_2 f_y = -\frac{1}{2} \log_2 f_X$

$\log_2 \left(\left(1 + \dfrac{M_y}{L}\right) 2^{E_y - B}\right) = -\frac{1}{2} \log_2 \left(\left(1 + \dfrac{M_X}{L}\right) 2^{E_X - B}\right)$
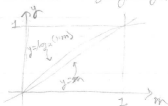
$\log_2 \left(1 + \dfrac{M_y}{L}\right) - \log_2 \left(2^{E_y - B}\right) = -\frac{1}{2}\left[\log_2\left(1 + \dfrac{M_X}{L}\right) + \log_2\left(2^{E_X - B}\right)\right]$

$\log_2 \left(1 + \dfrac{M_y}{L}\right) + E_y - B = -\frac{1}{2}\log_2\left(1 + \dfrac{M_X}{L}\right) - \frac{1}{2}E_X + \frac{1}{2}B$

$$\log_2\left(1+\tfrac{M_y}{L}\right) + E_y = -\tfrac{1}{2}\log_{02}\left(1+\tfrac{M_x^2}{L}\right) - \tfrac{1}{2}E_x + \tfrac{3}{2}\varepsilon$$

$$\log_{0n}\left(1+\tfrac{M_y^2}{L}\right) + E_y = \tfrac{1}{2}\log_2\left(1+\tfrac{M_x^2}{L}\right) - \tfrac{1}{2}(E_x - 3\varepsilon) \qquad (2)$$

Now consider the binary $\log_2(1+m)$ for $0 \le m < 1$:



We see that $\log_2(1+m) \approx m$ for $0 \le m < 1$

More specifically, $\log_2(1+m) = m + \theta_m$ for $0 \le m < 1$ where some small error term $\theta_m$ $\qquad (3)$

$\therefore$ Substituting (3) into (2), we have:

$$\tfrac{M_y}{L} + \theta_y + E_y = -\tfrac{1}{2}\left(\tfrac{M_x^2}{L} + \theta_x\right) - \tfrac{1}{2}(E_x - 3\varepsilon)$$

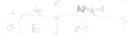$$\tfrac{M_y}{L} + E_y = -\tfrac{1}{2}\tfrac{M_x^2}{L} - \tfrac{1}{2}\theta_x - \theta_y - \tfrac{1}{2}(E_x - 3\varepsilon)$$

$$E_y L + M_y = -\tfrac{1}{2}M_x - L\left(\tfrac{1}{2}\theta_x + \theta_y\right) - \tfrac{L}{2}(E_x - 3\varepsilon)$$

$$= -\tfrac{1}{2}M_x - L\left(\tfrac{1}{2}\theta_x + \theta_y\right) - \tfrac{1}{2}E_x L + \tfrac{3}{2}\varepsilon L$$

$$E_y L + M_y = L\left(\tfrac{3}{2}\varepsilon - \left(\tfrac{1}{2}\theta_x + \theta_y\right)\right) - \tfrac{1}{2}(E_x L + M_x) \qquad (4)$$

Now lets take a quick look at an mbit floating point again again:



Point to an integer type, since $E2^{N-b-1} + M = EL + M$

So if we take $I_x$ and $I_y$ to be the floating point values $f_x + f_y$, respectively, cast to integers, then we have

$$I_x = E_x L + M_x \qquad , \qquad I_y = E_y L + M_y.$$  (6)

So from (5), (6) turns into

$$I_y = L\left(\tfrac{3}{2}B - \left(\tfrac{1}{2}\theta_x + \theta_y\right)\right) - \tfrac{1}{2}I_x$$

or more simply

$$I_y = R - \tfrac{1}{2}I_x$$

where $\quad R = L\left(\tfrac{3}{2}B - \left(\tfrac{1}{2}\theta_x + \theta_y\right)\right)$

And that's the basic technique: take floating point $x$ as an integer, divide it by 2 (right shift by 1), & subtract it from our magic integer $R$. Treating the result $I_y$ as a floating point, you get $y$.

The error comes from the fact that $\theta_x + \theta_y$ are for specific values of $x$ (& hence of $y$).

So we need to pick $\theta_x + \theta_y$ as best as we can, & it will be an approximation for most values of $x$.

For single precision IEEE floating point, $M = 23$, $B = 8$.

$$L = 2^{23}$$
$$B = 127$$

$$R = 2^{23}\left(\tfrac{3}{2}(127) - \left(\tfrac{1}{2}\theta_x + \theta_y\right)\right)$$